K ) Means Clustering        ⟶ whatsapp group
   ↳ (No of clusters)

⊕ Clusters ⟶ high Intercluster
           ⟶ low Intracluster distance    } ⟶ $DI = \dfrac{min(Int)}{max(Intra)}$

K Means Clustering

(Algo)
⟶ simple
⟶ K = no of clusters

Hyper parameter ⟹ which you give to the model

⊕  ⟋

$S_1 U S_2 U S_3 = Data(All)$
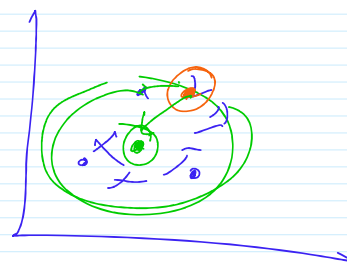$S_1 \cap S_2 \cap S_3 = \phi$

$C_3$  $S_3$

$C_1$  $S_1$     $C_2$  $S_2$

Centroid  (1) ⟶

| $x_1$ | $x_2$ | $x_3$ |
|-------|-------|-------|
| 2     | 3     | 5     |
| 3     | 8     | 9     |
| 4     | 9     | 10    |

(2) ⟶
(3) ⟶

Plot ⟶

$\dfrac{2+3+4}{3}$  $\dfrac{3+8+9}{3}$  $\dfrac{5+9+10}{3}$

↓

[ One point in a cluster can only belong to one cluster ]
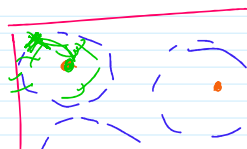
Goal ⟶ we have to find $C_1\, C_2 --- C_K$
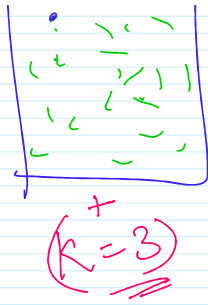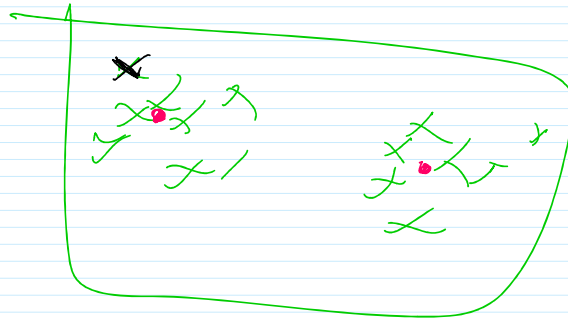         (K-clusters)

Algo

magic Box

K Means

(Assign each

Algo



K Means

(K = 3)

Mathematical Formulations
( k → we will give)

Goal: Find $c_1 \, c_2 \, c_3 \, ---\, c_K$

S.t. ① → Interclusters is maximised Distance

② minimise Intercluster distance

$S_1 \cup S_2 \cup S_3 = \emptyset$

$\widehat{S_i} \rightarrow \boxed{c_i}$

✓ $S_i \cap S_j = \phi$

each $x_i \in$ one $S_j$

$f_2$    $S_3$

$S_1$    $S_2$

$f_1$

Mean data pt. { $Centroid_1 = \dfrac{x_1 + x_2 + \cdots + x_m}{m} \in S_1$

Optimisation ⟶ (NP hard) → exponential time complexity
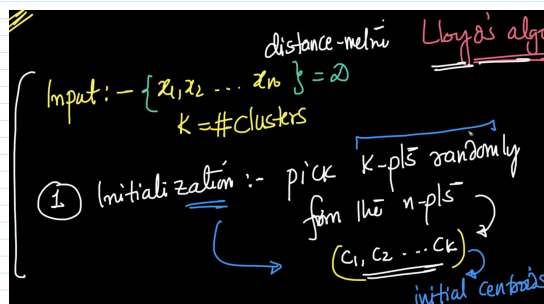
(#) Each data pt. is assigned only to one cluster

Lloyd's Algo
(Approximate)

(*) Lloyds Algorithm → only 3 step Process

① Random initialise points (k points)

② Assign each point to nearest cluster

③ Recompute / Update the centroid

distance-metric    Lloyd's algo

Input :- $\{x_1, x_2 \ldots x_n\} = \emptyset$

    K = #Clusters

① Initialization :- pick K-pts randomly

    from the n-pts

    $(c_1, c_2 \ldots c_k)$

    initial centroids

① Randomly initialise points

K

② Calculate distance of $\begin{bmatrix} x_i \\ i=1 \text{ to } N \end{bmatrix}$ from each centroid & then assign the label

② Assignment — Data — dist-metric

for each $x_i$ in ②
- Select the nearest centroid: $\tilde{c}_j$ (let)
- add $x_i$ to $S_j$
  ↳ assign that $x_i$ point to $S_j$ cluster

Update each of the $K$ cluster ⟹

③ Recompute / update the centroids → means ⟹ prone to outliers

for j in 1 to k
→ $c_j = \frac{1}{|S_j|} \sum_{x_i \in S_j} x_i$

mean-pt

---

K Means (Lloyds Algo) Scratch Implementation
(k=3)

① Random initialisn → n/p-random. choice (       )

② Dist ($x_i$, with each centroid)

$x_i \nearrow^{k_1}$ (Euclidian)
$\rightarrow k_2$ (min) ⟹ $x_i \longrightarrow k_2$
$\searrow k_3$

③ Recompute centroids $\left( \frac{\sum x}{n} \right) \Big/ \left( \frac{\sum y}{n} \right)$

$K \rightarrow$ hyperparameters (you have to give K)

How to find the right K ?

* (Business Sense)          * (Elbow Method)

↳ General concept in clustering : WCSS (within clusters sum of square)

Fashion tag its products

/↑ → ...

Car Rental company

↳ CarA (Hatchback)
(SUV)

1

products

Shirt  Jeans  Night Gown
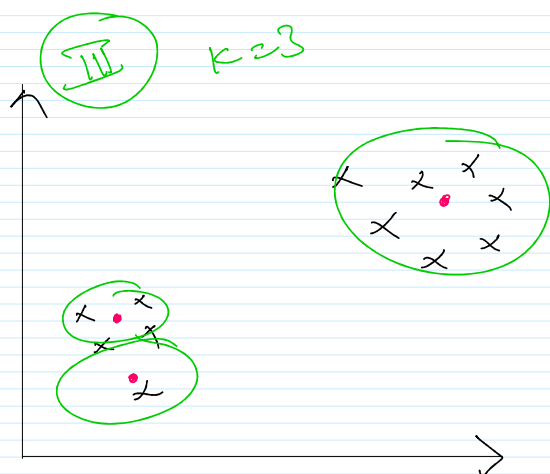
→ (K=3)

→ Car A (Hatchback)
→ Car B (SUV)
→ Car C (Luxury)

K = 3

$$WCSS = \sum_{P_i \text{ in Cluster 1}} distance(P_i, C_1)^2 + \sum_{P_i \text{ in Cluster 2}} distance(P_i, C_2)^2 + \sum_{P_i \text{ in Cluster 3}} distance(P_i, C_3)^2$$

(Effect of $K$ on WCSS)

no of clusters

I

$f2$

centroid

$f1$

II

$f2$

$C_1$

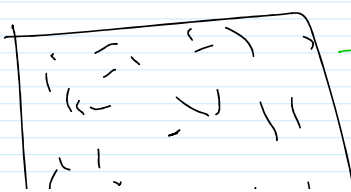$C_2$

$f_1$

(a)  WCSS  ① > ②

(b)  ① < ②

(c)  ① = ②

III  K=3

# [As $K \uparrow$ my WCSS decreases]

Quiz  when WCSS = 0 ?

① K = ?
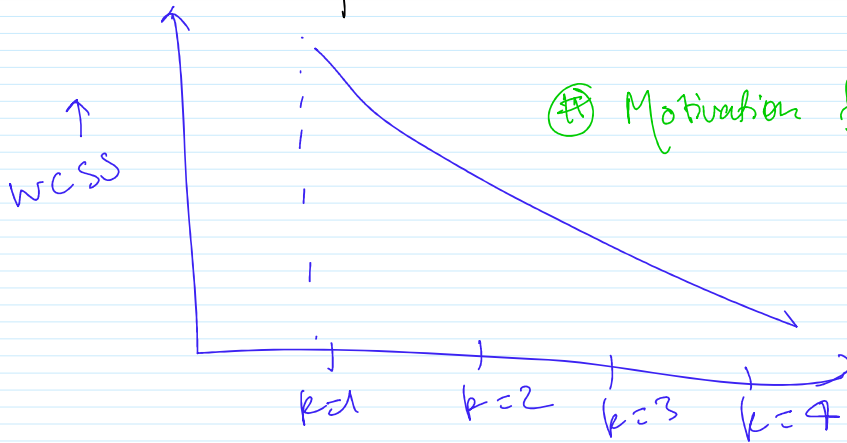
$n = (100$ data points$)$

Elbow Method  → WCSS  for $k = 1$ to say (5)

(3 — 4)
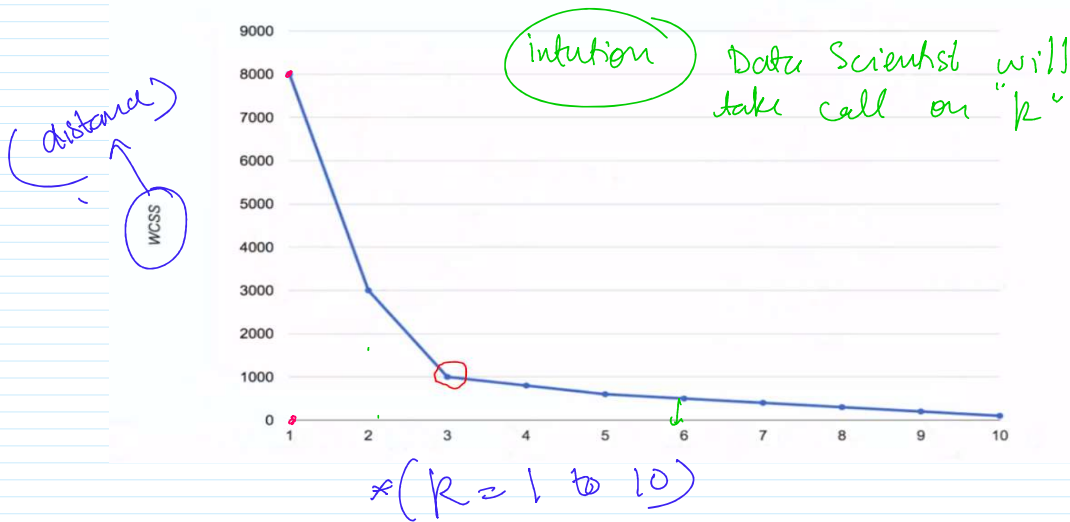
Unsupervides K means Page 4

$$(3 - 4)$$



WCSS

k=1  k=2  k=3  k=4

Motivation for finding good k ⇒
Slope at each point



(distance)

WCSS

intuition    Data Scientist will
            take call on "k"

$* (k = 1 \text{ to } 10)$



k means

(k = ? )

→ Business

Trap Points → { Edge Cases of Algorithm }

Low Dimension → Use Euclidean
Low — Med → Manhattan
— high → Cosine